

# A continuing overview of STAT 20 R

JE Hug

9/16/2020

Welcome to this document created by Josh Hug, one of the GSI's for STAT 20. This document will try and keep up with our R progress as we go on and I will add in depth (time permitting) explanations of what we are doing as I go on. Remember to check github for the latest updates to this document as it progresses!

A generic glossary of R commands that we have used so far.

```
# these are the packages we are using so far
```

```
library(dplyr)
library(ggplot2)
```

I'll begin with some simple use of the main dplyr functions we use, on the palmer penguins data set (make sure to install it if you don't have it yet). I prefer using this dataset over something like iris due to the fact that while iris is a classic dataset, it was compiled by Ronald Fisher (a fervent eugenicist) and published in a eugenics journal originally. This dataset provides a nice alternative with similar properites.

## 1 Filter, Select, Mutate Basics

The key facts here are to use select if we want

```
# Here I will use the palmer penguins data set
```

```
# if you don't have it installed
```

```
# install.packages("palmerpenguins")
```

```
library(palmerpenguins) # where this data set comes from
```

```
glimpse(penguins) # a nice function to take an easy look at the data
```

```
## Rows: 344
## Columns: 8
## $ species      <fct> Adelie, Adelie, Adelie, Adelie, Adelie, Adelie, A...
## $ island       <fct> Torgersen, Torgersen, Torgersen, Torgersen, Torge...
## $ bill_length_mm <dbl> 39.1, 39.5, 40.3, NA, 36.7, 39.3, 38.9, 39.2, 34....
## $ bill_depth_mm <dbl> 18.7, 17.4, 18.0, NA, 19.3, 20.6, 17.8, 19.6, 18....
## $ flipper_length_mm <int> 181, 186, 195, NA, 193, 190, 181, 195, 193, 190, ...
## $ body_mass_g   <int> 3750, 3800, 3250, NA, 3450, 3650, 3625, 4675, 347...
## $ sex          <fct> male, female, female, NA, female, male, female, m...
## $ year         <int> 2007, 2007, 2007, 2007, 2007, 2007, 2007, 2007, 2...
```

Suppose I wanted to make a new column bill\_length\_cm and body\_mass\_kg (where I convert units into cm and kg respectively)

We can use mutate to add a new column as some function of another column

```
new_pen <- mutate(penguins, bill_length_cm = bill_length_mm / 10, body_mass_kg = body_mass_g / 1000 )  
  
glimpse(new_pen)
```

```
## Rows: 344  
## Columns: 10  
## $ species      <fct> Adelie, Adelie, Adelie, Adelie, Adelie, Adelie, A...  
## $ island        <fct> Torgersen, Torgersen, Torgersen, Torgersen, Torge...  
## $ bill_length_mm <dbl> 39.1, 39.5, 40.3, NA, 36.7, 39.3, 38.9, 39.2, 34....  
## $ bill_depth_mm <dbl> 18.7, 17.4, 18.0, NA, 19.3, 20.6, 17.8, 19.6, 18....  
## $ flipper_length_mm <int> 181, 186, 195, NA, 193, 190, 181, 195, 193, 190, ...  
## $ body_mass_g    <int> 3750, 3800, 3250, NA, 3450, 3650, 3625, 4675, 347...  
## $ sex            <fct> male, female, female, NA, female, male, female, m...  
## $ year           <int> 2007, 2007, 2007, 2007, 2007, 2007, 2007, 2007, 2...  
## $ bill_length_cm <dbl> 3.91, 3.95, 4.03, NA, 3.67, 3.93, 3.89, 3.92, 3.4...  
## $ body_mass_kg   <dbl> 3.750, 3.800, 3.250, NA, 3.450, 3.650, 3.625, 4.6...
```

Now suppose I want penguins that weigh less than like 3000 g (3kg) only. Since this is subsetting over rows with a specific condition we use the filter function from dplyr.

```
new_pen_light<- filter(new_pen, body_mass_kg <3 )  
  
glimpse(new_pen_light)
```

```
## Rows: 9  
## Columns: 10  
## $ species      <fct> Adelie, Adelie, Adelie, Adelie, Adelie, Adelie, A...  
## $ island        <fct> Dream, Biscoe, Biscoe, Biscoe, Dream, Biscoe, Tor...  
## $ bill_length_mm <dbl> 37.5, 34.5, 36.5, 36.4, 33.1, 37.9, 38.6, 43.2, 46.9  
## $ bill_depth_mm <dbl> 18.9, 18.1, 16.6, 17.1, 16.1, 18.6, 17.0, 16.6, 16.6  
## $ flipper_length_mm <int> 179, 187, 181, 184, 178, 193, 188, 187, 192  
## $ body_mass_g    <int> 2975, 2900, 2850, 2850, 2900, 2925, 2900, 2900, 2700  
## $ sex            <fct> NA, female, female, female, female, female, femal...  
## $ year           <int> 2007, 2008, 2008, 2008, 2008, 2009, 2009, 2007, 2008  
## $ bill_length_cm <dbl> 3.75, 3.45, 3.65, 3.64, 3.31, 3.79, 3.86, 4.32, 4.69  
## $ body_mass_kg   <dbl> 2.975, 2.900, 2.850, 2.850, 2.900, 2.925, 2.900, ...
```

## 2 Histograms in R

Take a look at this cheat sheet, you're probably extremely overwhelmed by this and personally I don't know what at least half of the stuff on this page does but it will save you a lot of time from googling. There are actually cheat sheets for most tidyverse (what these packages are a part of) packages so you can check out ones for dplyr and such.

### 2.1 General format of ggplot

The general format of ggplot is that we call some generic function, then we can just build on it by (literally) adding to it other components. Now we always start with the base ggplot function which has two main inputs the dataframe and then the aesthetic. ggplot is extremely flexible and I can't begin to scratch the surface of what you can do here so take a look at some other examples or the cheat sheet.

Suppose I want a histogram of penguins bill length in mm. I tell ggplot to look at my dataframe penguins and then I tell it to look for a specific column by its column name. Note that I don't have to subset my

dataframe at all before this if I want to use entire columns, no matter how many other columns ggplot only looks at the ones I tell it to.

```
ggplot(penguins, aes(x=bill_length_mm))
```

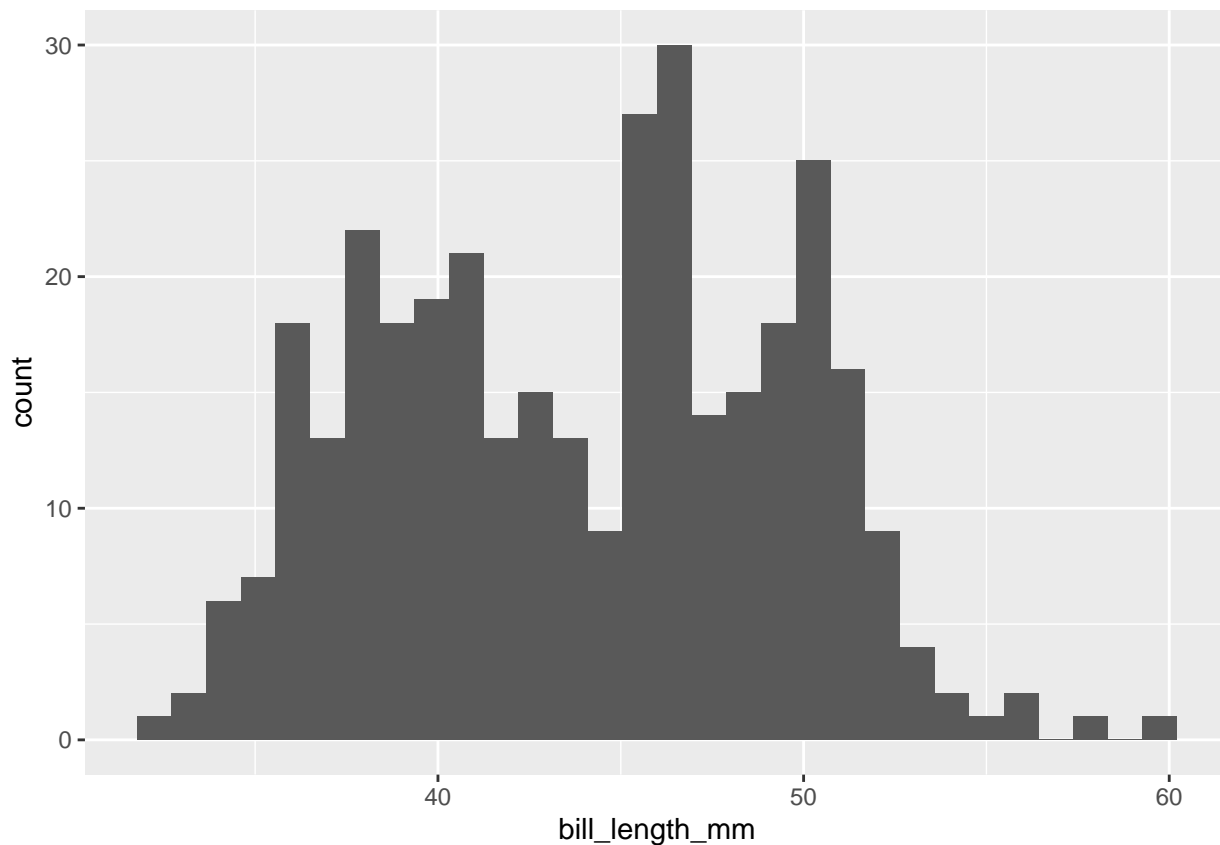


So R here has outputted nothing because I have told it only what the data frame and the column I want inputted is but I haven't specified which type of plot I want. So my next line I literally add onto it to tell R to make a histogram

```
ggplot(penguins, aes(x=bill_length_mm)) +  
  geom_histogram() # this is the line that tells it to make a histogram
```

```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```

```
## Warning: Removed 2 rows containing non-finite values (stat_bin).
```

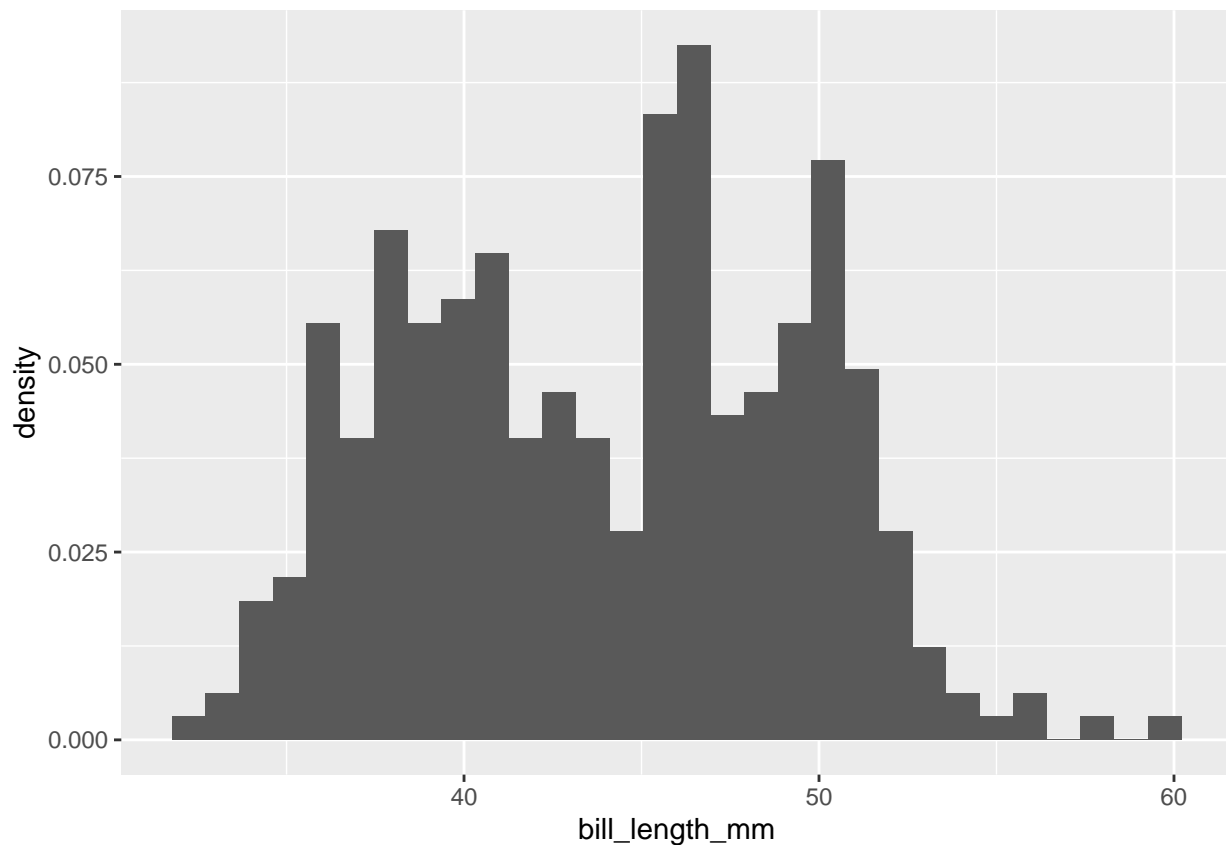


So R has made a histogram here but notice that we don't have density on the y axis we just have counts which is not what we usually want, so we can add another argument to aes that specifically tells it to use the density on the y axis.

```
ggplot(penguins, aes(x=bill_length_mm, y=..density..)) +  
  geom_histogram() # this is the line that tells it to make a histogram
```

```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```

```
## Warning: Removed 2 rows containing non-finite values (stat_bin).
```

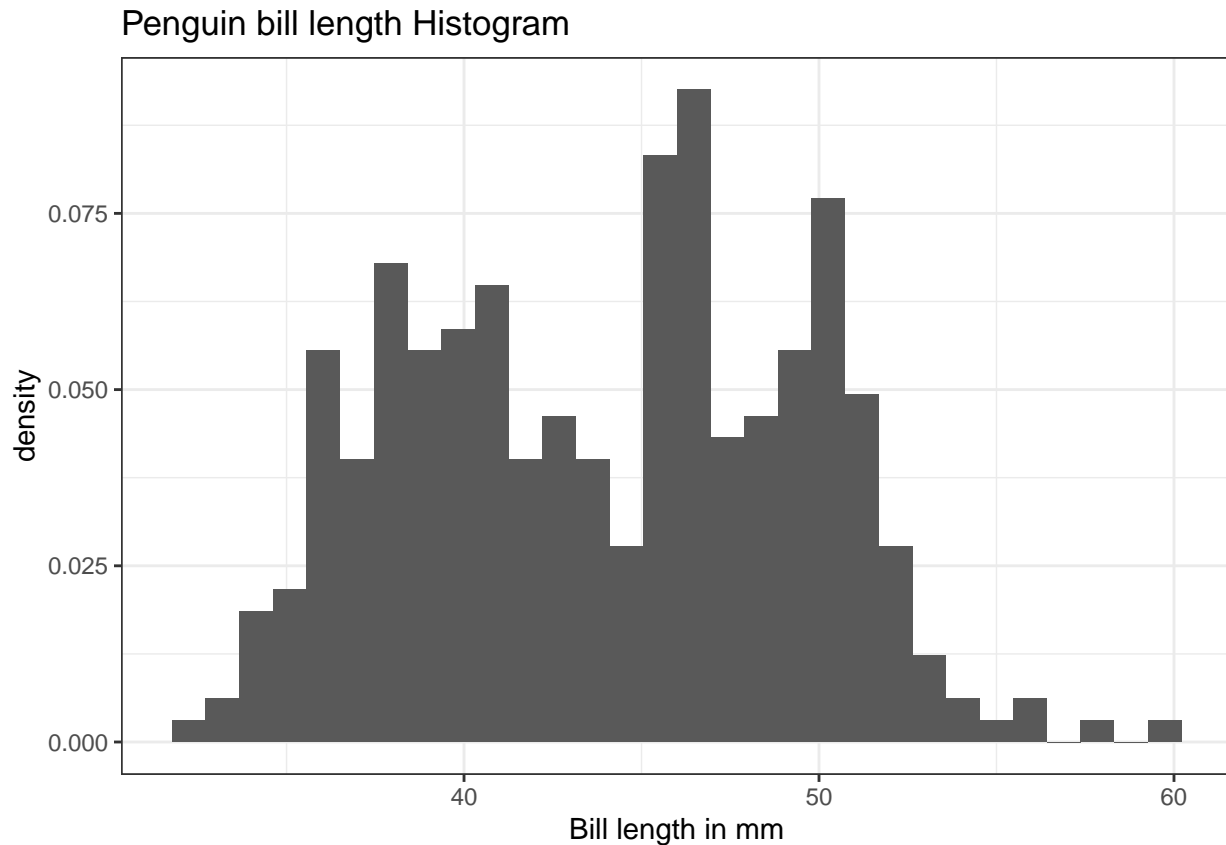


Now we have a histogram like we want it. I will add a title, change the theme and change the axis titles and we can be done for now it's really that simple.

```
ggplot(penguins, aes(x=bill_length_mm, y=..density..)) +
  geom_histogram() + # this is the line that tells it to make a histogram
  xlab("Bill length in mm") + # changing the x axis label
  ggtitle("Penguin bill length Histogram") + # adding the title
  theme_bw() # changing the theme for fun
```

```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```

```
## Warning: Removed 2 rows containing non-finite values (stat_bin).
```



## 3 The Rep function and sampling

### 3.1 A simple sampling example

Now R is made for statistics and it has amazing random sampling functionality. If you need to sample from something oftentimes you can. Sampling is very important in any statistical practice since it allows us to generate fake data. When I say sampling I essentially mean simulating data. For instance I can simulate a coin toss many times. Let 1 denote heads and 0 denote tails.

```
coin <- c(0,1)
```

From now on for reproducibility I will set a seed (what this does is allows for you to get the exact same random results as I get when you run these functions).

```
set.seed(11)
```

My coin here is simply a vector with a 0 and a 1. I now need a function to randomly pick one of those with equal probability, This is the sample function

```
sample(coin ,size=1)
```

```
## [1] 1
```

So I have drawn heads here, but this is slow for me to run this manually so many times! This is what the size option is for. Let's say I wanted to flip 100 coins.

```
sample(coin, size=100)
```

```
## Error in sample.int(length(x), size, replace, prob): cannot take a sample larger than the population
```

We've received an error that says "cannot take a sample larger than the population when `replace = FALSE`". This is because what R does by default is that it treats your vector like a hat we're pulling papers out of. First it goes in and pulls out a random entry from the vector (in our case this was the 1), but by default it doesn't put this back into the hat, so then the next draw from the vector hat is whatever is left which in our case is just the 0. Finally on the third draw there are no numbers left in the hat to draw from so R spits out an error. We want R to put the papers back into the hat as if we are flipping the coin brand new. We can do this with the `replace=TRUE` argument.

```
flips<-sample(coin, size=100, replace=TRUE)
flips
```

```
## [1] 1 1 0 1 0 0 1 1 1 0 1 1 1 0 1 0 0 0 0 0 1 0 0 1 0 0 1 1 0 0 0 1 1 0 1 0
## [38] 1 1 1 0 1 0 1 1 0 1 1 0 1 0 1 0 0 0 1 1 0 1 1 0 0 1 1 1 1 0 1 0 0 0 1 0 1
## [75] 1 1 0 1 1 1 1 0 1 0 0 1 0 1 0 1 1 0 0 0 1 0 0 0 1 0
```

So we've now received our coin draws. Hopefully since we know this is a fair coin (why?), we can assume that we will have about 50 heads. I can just sum to check

```
sum(flips)
```

```
## [1] 51
```

That's pretty close to 50. In addition we can see the proportion of heads by dividing by the length.

```
sum(flips)/length(flips)
```

```
## [1] 0.51
```

An important thing to note is that this is exactly the same as taking the mean which we define as the sum of all our observations divided by the length

```
mean(flips)
```

```
## [1] 0.51
```

### 3.2 More complicated sampling with rep

Suppose I wanted to flip a coin with probability 1/56 of landing on heads. We can make this with a vector of 55 zeros and a single one. Note that in defining `weird_coin` I use `rep`, because I'm putting the same value so many times I don't want to write it out so I tell with `rep(0,55)` give me a vector with 55 zeros then with `c(rep(0,55),1)` I append a single one to the new vector.

```
weird_coin <- c(rep(0,55),1)
```

So here I calculate the mean of this (again the same as the proportion as shown above), this should be fairly close to 1/56

```
close<-mean(sample(weird_coin,1000,replace=TRUE))
```

```
close-(1/56)
```

```
## [1] -0.005857143
```

It is pretty close.

Someone asked about flipping tails rather than heads well in this case if I made the probability of heads 1, I can just do  $1 - P(\text{heads})$  (the complement rule) or I could replace the vector with the complementary vector. This is just we switch the ones and zeroes.

```
weird_coin_tails <- 1-weird_coin # Why does this work?
```

```
close_tails<-mean(sample(weird_coin_tails,1000,replace=TRUE))
```

```
close_tails-(1-(1/56))
```

```
## [1] 0.003857143
```

Again this is pretty close to what we'd expect it to be.

### 3.3 an extension of what we know (simulating things we don't know the answer to)

Sampling is a rather simple function but we can sometimes make it cater to some more complicated problems however almost all of them can essentially be boiled down to coin flips of some kind. For instance on a previous homework we were mostly doing simulations to confirm something we already knew. Here is an example of simulating with sampling to show something more complicated. All of this can be solved by hand but would be tedious. Don't worry if you don't understand the coding here (in fact its coded rather poorly by me) but this is trying to see the probability of rolling three die and seeing the probability of getting a dice roll 10 or higher. Inside of the for loop is what is important, I'm sampling three draws with replacement from our dice then summing them together. I am then adding these sums to a vector. After this I take the mean of the number of times the sum was greater than 10. This should approximate the probability of rolling 3 dice and having a sum greater than 10. (this is mathematically supported by something called the law of large numbers).

```
dice<- c(1:6)
```

```
z <- replicate(100000,sum(sample(dice,replace=TRUE,size=3)))
```

```
mean(z>=10)
```

```
## [1] 0.62501
```

## 4 Probability

### 4.1 A brief introduction

I will assume knowledge of basic probability rules learned in lecture such as multiplication rule and addition rule. Key among these is the definition of independence which I will emphasize here. Denote

$$P(A \text{ and } B) \equiv P(A, B)$$

. Recall that two events  $A, B$  are independent if and only if

$$P(A, B) = P(A)P(B)$$

Suppose that  $P(A) \neq 0$ . We know by the multiplication rule that

$$P(A, B) = P(B|A)P(A)$$

Combining these two facts we see that

$$\begin{aligned} P(B|A)P(A) &= P(B)P(A) \\ P(B|A) &= P(B) \end{aligned}$$

which can be interpreted as the fact that knowing  $A$  has no effect on the probability of  $B$  happening!



## 4.2 Binomial Distribution

The easiest way to think of the binomial distribution is that it is simply a sum of independent coin flips. Problems solved with this distribution are of the form “What is the probability of getting 6 heads in 10 coin flips”. Let  $X$  be the sum of the results of  $n$  coin flips each with probability of heads being  $p$ . Then we say  $X \sim \text{binomial}(n, p)$  which is read that  $X$  has the probability distribution binomial with parameters  $n$  and  $p$ . We can then compute the probability that  $X = k$  (the probability of  $k$  heads in  $n$  tosses with each toss being independently heads with probability  $p$ ), written  $P(X = k)$ . We have that

$$P(X = k) = \binom{n}{k} p^k (1 - p)^{n-k}$$

Note that the  $p^k(1 - p)^{n-k}$  is simply the probability of having the heads show up in one specific spot within the  $n$  tosses. The  $\binom{n}{k}$  accounts for the other spots we can put the successful tosses. Therefore the binomial distribution is both an application of multiplication and the addition rule with both independence and mutually exclusive events (Don't worry if you don't understand this).

## 4.3 Computing binomial probabilities in R

The `dbinom` function in R evaluates the binomial formula above it takes in arguments `x`, `size` and `prob`. For our purposes `x` refers to  $k$ , `size` refers to  $n$  and `prob` refers to  $p$ .

So if I wanted the probability of seeing 5 heads in 10 fair coin tosses I would evaluate

```
dbinom(x=5,size=10,prob=1/2)
```

```
## [1] 0.2460938
```

The `choose` function on the other hand simply evaluates  $\binom{n}{k} = \frac{n!}{(n-k)!k!}$  (The function takes arguments `n` and `k` as in the formula). So if I wanted to evaluate the same probability I could do:

```
# the same probability as above
```

```
choose(10,5) * (0.5)^(5) * (1-0.5)^(10-5)
```

```
## [1] 0.2460938
```

This is clearly more tedious and the `choose` function should really only be used if you need to use `choose` specifically.

### An interlude on the binomial coefficient

$\binom{n}{k} \equiv \frac{n!}{(n-k)!k!}$  read “ $n$  choose  $k$ ” is a special value in combinatorics which is the number of ways to choose an unordered subset of  $k$  elements from a set of size  $n$ . For instance the number of teams of 3 people I can make from 9 people total is  $\binom{9}{3} = 84$  (check with R!).

## 4.4 The difference between `rbinom`, `dbinom` and `pbinom`

This is copy pasted from a piazza response I wrote recently so hopefully this helps!

`dbinom(x=k,size=s,prob=p)` literally just evaluates the binomial formula  $\binom{s}{k} p^k (1 - p)^{s-k}$ . So we use this when we want the probability of some binomial random variable. If  $X$  is a binomial random variable with probability of success  $p$  and  $s$  number of trials, then the probability of  $k$  successes is given by this or  $P(X = k)$ .

`pbinom(q=k,size=s,prob=p)` evaluates a sum of binomial probabilities ( a sum of `dbinom`) for all values **less than or equal** to  $k$ . So in this case `pbinom(k,s,p)=sum(dbinom(c(0:k),s,p))`. So similar set up to `dbinom` we have that this gives the probability of having anywhere from 0 to  $k$  successes, out of  $s$  trials each with probability of success  $p$ , or  $P(X \leq k)$ .

The two above functions have evaluated the probability of seeing certain events, however `rbinom` is used to generate events or sample. Essentially R is flipping coins for us and then giving us the results. Take the simplest example `rbinom( n=1 , size=1, prob=1/2)` this will return either 1 or 0 with probability  $\frac{1}{2}$ , this is the result of a single coin flip with heads returning 1 and tails returning 0. If we want a sum of say  $s$  coin flips we can change this to `rbinom(n=1,size=s,prob=1/2)`. This returns a number between 0 and  $s$  representing the number of times the coin landed heads. Finally suppose I wanted to run  $z$  experiments where in each experiment I flip  $s$  coins, this is given by `rbinom(n=z, size=s, prob= 1/2)` where this returns a vector of length  $z$  where each entry is the number of heads in flipping  $s$  coins. The only parameter I haven't mentioned is `prob` which is just the probability of seeing a heads (a 1). Therefore we can change this to anything we want and have `rbinom(n=z, size=s, prob= p)` which is telling R, to flip  $s$  coins each with probability of heads  $p$  record the number of heads, then repeat this  $z$  times and return a vector with the results.

## 5 R Markdown

R Markdown is one of my favorite tools. If you've been to my section or to my office hours you've seen me use it since I don't usually code within R scripts. In fact this document is written in R markdown. Working within chunks allows for easily organized code and reproducibility. In addition all the math here is embedded via LaTeX, a very popular way to type complicated math. As usual I recommend never working in the console except for specific instances of execution of R code. The best way to learn R Markdown (like R in general) is to just play around with it and figure out what works. You probably need a local LaTeX installation to knit to pdf, but you all should be able to knit to html without that.

## 6 The Law of Averages (or why I prefer the law of large numbers)

If you were to google the Law of Averages you'd note that in actuality it is not a mathematical statement. My background is more mathematical statistics therefore in general I try to be more accurate with you all about actual statements about probability. For full context I've never been taught the law of averages, so I will try and convert it into language that I find more familiar and more accurate.

As stated in class we have generally that for a generic binary (the outcome is a success or failure) we have that

$$1) \# \text{ successes} = \text{Expected } \# \text{ of successes} + \text{Chance Error}$$

The chance error is kind of the "randomness" for instance I'll flip fair coin 100 times and print out the chance error

```
flips <- rbinom(1,100,0.5)
flips-50
```

```
## [1] -4
```

So our chance error is what is returned. Now I'm going to make a statement about the general trend of this chance error, but we always need to be careful about these statements because in reality (mathematically) all of these statements only actually hold asymptotically (when we are at infinity). As we increase the sample size (aka the number of trials here) we should **expect** a larger absolute chance error. Now this doesn't mean that we will see it, it's entirely possible not but we should expect to.

```
flips2 <- rbinom(1,10000,0.5)
flips2-5000
```

```
## [1] -1
```

So we've seen that the chance error is bigger like we'd expect. This should make sense since we are doing a lot of trials there is potential for things to go wrong.

$$2) \% \text{ successes} = \text{Expected } \% \text{ of successes} + \% \text{ chance error}$$

This is the much more meaningful in my opinion because it is equivalent to, letting  $n$  denote the sample size

$$2) \frac{\# \text{ Successes}}{n} = \frac{\text{Expected } \# \text{ of successes}}{n} + \frac{\text{Chance error}}{n}$$

One should note that all three numerator terms are functions of  $n$ .

But as we take  $n$  to be larger we should **expect** to get smaller, since the probability of being far from the expected  $\#$  of successes becomes smaller. If the sample size goes to infinity this is an application of the law of large numbers, which states that we expect any sample average to converge to the actual expected value (with some regularity conditions).

Note that as long as we can group into failures or successes, this can work so even if we have 3 outcomes and we look at one outcome as a success then we can adapt with this binary model, as everything but a subset to be a failure and the subset is the success effectively creating a binary problem.

## 7 Random Variables

I've had many people confused about random variables and what exactly they are. The technical definition of random variables is a bit tricky. A random variable is colloquially what the name entails, a variable that is random. Generally we denote random variables by uppercase letters like  $X, Y, Z$ . Now that is the actual random variable, we also denote instantiations, essentially what we observe by the lowercase counterpart. So generally  $X$  is the random variable the observed value is  $x$ . Random variables are always numbers, this is important because it allows us to do math on them. We've already seen a simple random variable, a coin toss in which heads is one and tails is 0. We can see quantify all possibilities of the random variable and its probabilities with a distribution table.

$x$	$P(X = x)$
0	0.5
1	0.5

So the random variable  $X$  takes on values  $x$  with probabilities as given in the table. Note that this can be represented as a function

$$P(X = x) = 0.5^x(1 - 0.5)^{(1-x)}$$

More generally for a weighted coin toss this becomes

$x$	$P(X = x)$
0	$1 - p$
1	$p$

which can also be represented as a function

$$P(X = x) = p^x(1 - p)^{(1-x)}$$

Hopefully this is familiar to you as a special case of the binomial formula, in fact it is the case where  $n = 1$  since the binomial is a sum of independent coin flips. Therefore if  $Y$  is binomial it has a probability distribution which can be represented by

$$P(X = x) = \binom{n}{x} p^x (1 - p)^{1-x}$$

Besides that most of the problems we have worked with so far can be expressed as random variables, in particular box models. A lot of questions we have focus on translating box models to random variables. The general methodology is to model a random variable  $X$  as one draw from the box model, from which you can usually easily compute probabilities. As will be shown in the next section we can fairly easily evaluate expected values (standard errors) of sums of (independent) random variables.

## 8 Expected Value and Standard error

When dealing with random variables (and box models), we know that we have distributions. These completely characterize the distribution, however they are a lot of information and we want to be able to know some elements of the distribution without these. Therefore oftentimes we are able to use summary statistics in order to better understand random variables. For now we are only really considering discrete random variables (these variables take on a finite<sup>1</sup> number of values) rather than continuous random variables that take on any real number in a specific interval. The binomial for instance can only take on values from 0 to  $n$ , by contrast the normal distribution random variables can take on any value on the real line. The definitions are given below, consider some discrete random variable  $X$  that takes on values of  $\mathcal{X}$ , then we have that

$$\mathbb{E}(X) = \sum_{x \in \mathcal{X}} x \cdot P(X = x)$$

This sum can be infinite, and just sums over all values taken on by  $X$ . The standard error is given by

$$SE(X) = \sqrt{\sum_{x \in \mathcal{X}} (x - \mathbb{E}(X))^2 \cdot P(X = x)}$$

These probably seem rather abstract, but can be interpreted in the exact same way as average and SD can be. The expected value is our best guess as to where the random variable will be at any point, essentially an approximation of the constant. The standard error is how far we expect the random variable to be from its own mean approximated as a constant. Random variables that follow known distributions typically<sup>2</sup> have simple formulas for their expectation and standard error. Typically we use the variance or  $SE^2$  for convenience since the square root is rather painful to work with mathematically. For the binomial distribution we have computed in lecture that if  $X \sim \text{bin}(n, p)$  then

$$\mathbb{E}(X) = np, \quad SE(X) = \sqrt{np(1-p)}$$

which aligns exactly with the box model representation of the box model. This is because a box model is a special case of the random variable. Can you think about how to represent the binomial distribution (sum of independent coin flips), as a box model?

## 9 The normal approximation (the central limit theorem)

The Central Limit Theorem is one of the most important probability theorems for the practice of statistics. It essentially guarantees that under certain conditions we can make a decent approximation for any data that we have in the real world. In lecture we spoke about sums of draws from box models as well as the binomial distribution being able to be approximated by the normal distribution. These are again special cases of the actual central limit theorem. Suppose we have many random variables, let all of these have the same distribution and be independent (in the binomial case you can think of each of these as a coin flip). Denote each of these random variables  $X_i$ , and suppose that  $E(X_i) = \mu$  and  $SE(X_i) = \sigma < \infty$ , since each of these have the same distribution they all have the same expected value and SE. Then the<sup>3</sup> central limit

---

<sup>1</sup>This should really say countable

<sup>2</sup>some distributions have infinite or undefined moments

<sup>3</sup>well one version

theorem states that if we let  $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$  then

$$\frac{(\bar{X} - \mu)}{\sigma/\sqrt{n}}$$

is approximately standard normal distributed. As  $n$  gets larger this is a better and better approximation. Note that this uses the fact that the average of random variables  $\bar{X}$  has expected value  $\mu$  and standard error  $\sigma/\sqrt{n}$ . The reason why this is useful is that we use sums of random variables all the time, and we know very well how to work with sums of random variables.

## 10.1 Scatter plots and linear regression

Sorry this is so delayed but this will be rather brief, but this post was mainly focused on R anyway. First some notes on linear regression in R. We can obtain a regression in R with the syntax `lm(y~x)`, where we wish to predict  $y$  from  $x$ . Here `lm` stands for linear model. This is an example from my section where I emphasize that the lines for predicting each variable from the other are different. So essentially we cannot use either line to go “backwards”. Note both lines intersect at the mean of both variables since that point is on both regression lines.

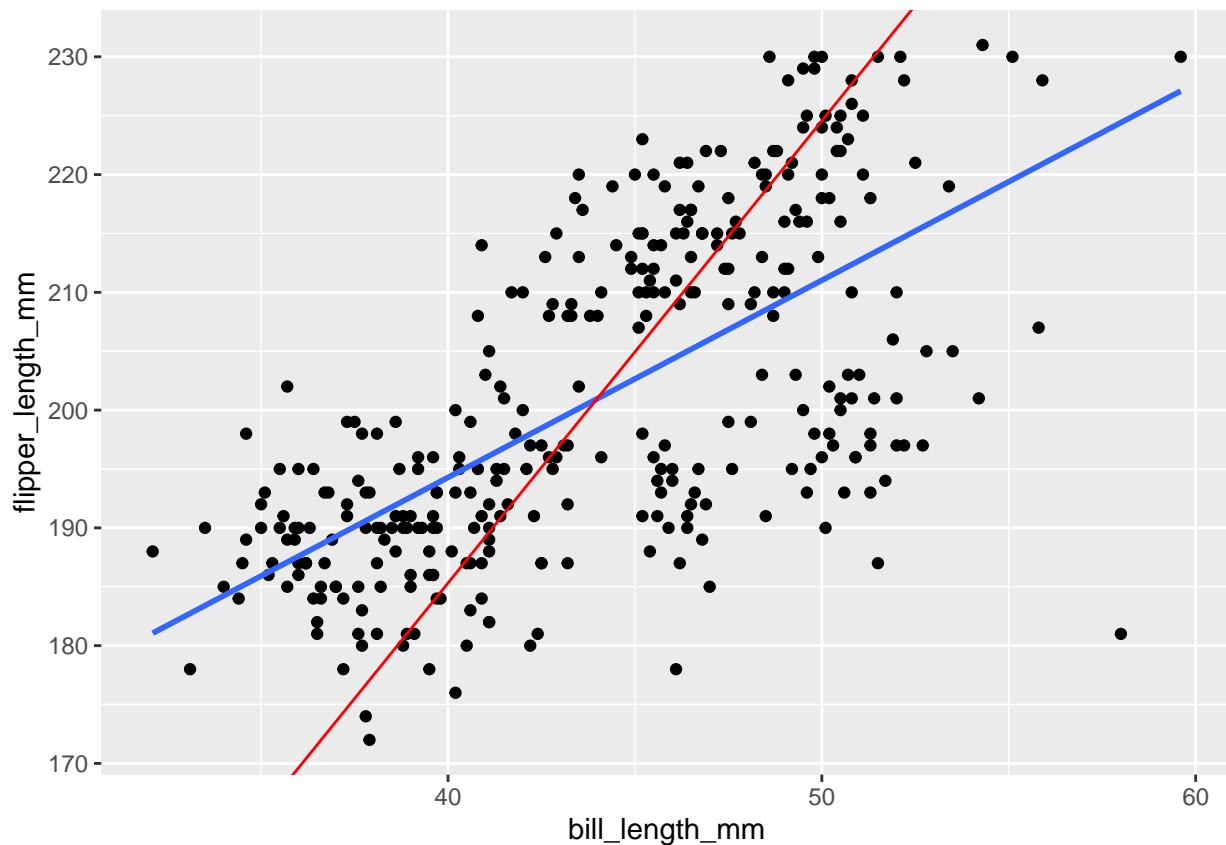
```
library(palmerpenguins)
library(ggplot2)

dat <- na.omit(penguins)

lin_mod <- lm(bill_length_mm~flipper_length_mm,data=dat)

ggplot(dat, aes(x=bill_length_mm,y=flipper_length_mm))+
  geom_point()+
  geom_smooth(method="lm", se= F)+
  geom_abline(intercept=7.2185580/0.2548, slope=1/0.2548, color="red")

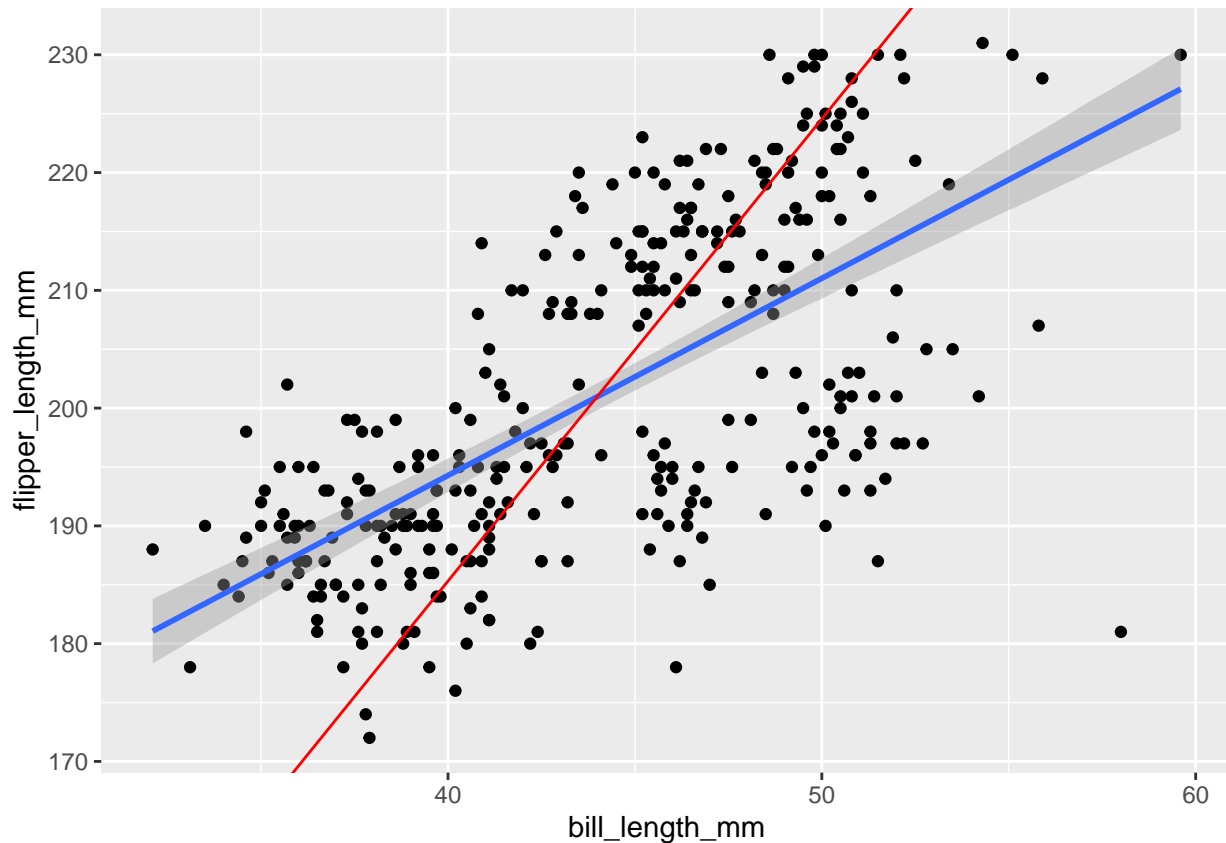
## `geom_smooth()` using formula 'y ~ x'
```



Building on our earlier ggplot syntax, we are just telling R to do a scatter plot with `geom_point`. Then `geom_smooth(method = "lm", se = F)` tells ggplot to plot a line based on the linear model and not to show standard error bars. If we omit the standard error option we get

```
ggplot(dat, aes(x=bill_length_mm,y=flipper_length_mm))+  
  geom_point()+  
  geom_smooth(method="lm")+ # tnote omitted SE=FALSE here  
  geom_abline(intercept=7.2185580/0.2548, slope=1/0.2548, color="red")
```

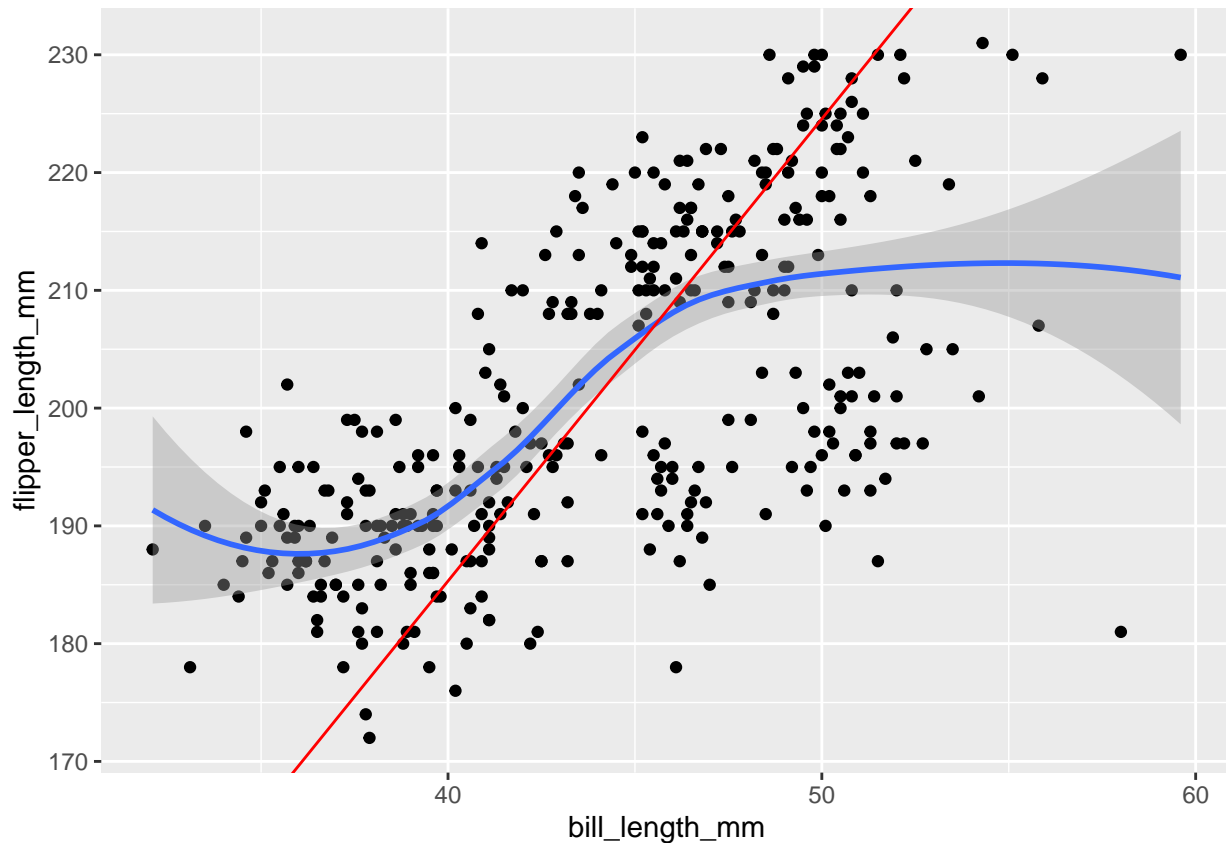
```
## `geom_smooth()` using formula 'y ~ x'
```



Where the grey bars are actually SE from the line. If we omit the `method="lm"` call ggplot by default will fit a non straight line with a method known as LOESS, much beyond the scope of this course but good to know.

```
ggplot(dat, aes(x=bill_length_mm,y=flipper_length_mm))+
  geom_point()+
  geom_smooth()+ # note omitted SE=FALSE here, and method = "lm"
  geom_abline(intercept=7.2185580/0.2548, slope=1/0.2548, color="red")
```

```
## `geom_smooth()` using method = 'loess' and formula 'y ~ x'
```



Clearly not a straight line.

## 10.2 How do I get the regression line in R?

This is a very simple call to the `lm` function for this.

```
# I omit the missing data (bad practice for actual analysis)
```

```
dat <- na.omit(penguins)
```

Suppose I wish to predict flipper length from bill length, the code for this regression line is simple.

```
myfit <- lm(formula = flipper_length_mm ~ bill_length_mm, data = dat)
```

```
myfit
```

```
##
## Call:
## lm(formula = flipper_length_mm ~ bill_length_mm, data = dat)
##
## Coefficients:
##   (Intercept)  bill_length_mm
##         127.330          1.674
```

This just returns the intercept and the slope of the line, pretty simple. For more complicated output we can take a look at the summary call of this regression.

```
summary(myfit)
```



```
##
## Call:
## lm(formula = flipper_length_mm ~ bill_length_mm, data = dat)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -43.413  -7.837   0.652   8.360  21.321
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   127.3304     4.7291   26.93  <2e-16 ***
## bill_length_mm  1.6738     0.1067   15.69  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 10.63 on 331 degrees of freedom
## Multiple R-squared:  0.4265, Adjusted R-squared:  0.4248
## F-statistic: 246.2 on 1 and 331 DF,  p-value: < 2.2e-16
```

This is a lot more complicated, and has some info that may confuse you. In the coefficients section we have the estimates column, these are the values for the slope and intercept of the line that we saw before. The next 3 columns are important for hypothesis testing. The hypothesis test we are performing here is

$H_0$ : Slope = 0

$H_1$ : Slope  $\neq$  0

as well as

$H_0$ : intercept = 0

$H_1$ : intercept  $\neq$  0

Note these are what we call two sided tests, not what we have done so far, but the interpretation is the same. Why are we interested in these tests? Well we want to see if we have a slope that is actually significantly different from 0, because if not there is no use for the line in predicting one from the other. The second column is the Standard error, and the third column is the t statistic value. We know how to calculate both of these. The final column is the p-value, here both are very small so we reject the null that the slope and intercept are 0. There is more info later on in the summary but it is most likely a bit beyond the scope of this course. In this simple case however the “Multiple R-squared” is just the correlation between the two variables squared.